

# 1Stepwise large genome assembly approach: A case of Siberian larch 2(*Larix sibirica* Ledeb.)

3

4Dmitry A. Kuzmin<sup>1,2</sup>, Sergey I. Feranchuk<sup>1,3,4</sup>, Vadim V. Sharov<sup>1,2</sup>, Alexander  
5N. Cybin<sup>1,2</sup>, Stepan V. Makolov<sup>1,2</sup>, Yuliya A. Putintseva<sup>1</sup>, Natalya V. Oreshko-  
6va<sup>1,5</sup>, Konstantin V. Krutovsky<sup>1,6,7,8\*</sup>

7

8<sup>1</sup>Laboratory of Forest Genomics, Genome Research and Education Center, Siberian Federal  
9University, 660036 Krasnoyarsk, Russia

10<sup>2</sup>Department of High Performance Computing, Institute of Space and Information Technologies,  
11Siberian Federal University, 660074 Krasnoyarsk, Russia

12<sup>3</sup>Department of Informatics, National Research Technical University, 664074 Irkutsk, Russia

13<sup>4</sup>Limnological Institute, Siberian Branch of Russian Academy of Sciences, 664033 Irkutsk,  
14Russia

15<sup>5</sup>Laboratory of Forest Genetics and Selection, V. N. Sukachev Institute of Forest, Siberian  
16Branch of Russian Academy of Sciences, 660036 Krasnoyarsk, Russia

17<sup>6</sup>Department of Forest Genetics and Forest Tree Breeding, Georg-August University of  
18Göttingen, 37077 Göttingen, Germany

19<sup>7</sup>Laboratory of Population Genetics, N. I. Vavilov Institute of General Genetics, Russian  
20Academy of Sciences, Moscow 119333, Russia

21<sup>8</sup>Department of Ecosystem Science and Management, Texas A&M University, College Station,  
22TX 77843-2138, USA

23

24\*Corresponding author: kkrutovsky@gmail.com; konstantin.krutovsky@forst.uni-goettingen.de

25

26

27

## 28Abstract

29**Background:** *De novo* assembling of large genomes, such as in conifers (~12-30 Gbp) that also  
 30consist of ~80% of repetitive DNA is a very complex and computationally intense endeavor. One  
 31of the main problems in assembling such genomes lays in computing limitations of nucleotide  
 32sequence assembly programs (DNA assemblers). As a rule, modern assemblers are usually  
 33designed to assemble genomes with a length not exceeding the length of the human genome  
 34(3.24 Gbp). Most assemblers cannot handle the amount of input sequence data required to  
 35provide the sufficient coverage needed for a high-quality assembly.

36**Results:** An original stepwise method of *de novo* assembly by parts (sets), which allows to  
 37bypass the limitations of modern assemblers associated with a huge amount of data being  
 38processed is presented in this paper. The results of numerical assembling experiments conducted  
 39using the model plant *Arabidopsis thaliana*, *Prunus persica* (peach) and four most popular  
 40assemblers, ABySS, SOAPdenovo, SPAdes and CLC Assembly Cell, showed the validity and  
 41effectiveness of the proposed stepwise assembling method.

42**Conclusion:** Using the new stepwise *de novo* assembling method presented in the paper the  
 43genome of Siberian larch, *Larix sibirica* Ledeb. (12.34 Gbp) was completely assembled *de novo*  
 44by the CLC Assembly Cell assembler. It is the first genome assembly for larch species in  
 45addition to only five other conifer genomes sequenced and assembled for *Picea abies*, *Picea*  
 46*glauca*, *Pinus taeda*, *Pinus lambertiana* and *Pseudotsuga menziesii* var. *menziesii*.

47**Keywords:** *de novo* genome assembly, Siberian larch, *Larix sibirica*

48

## 49Background

50The *de novo* assembling of large genomes, such as in conifers, that have a length of 12 to 30 Gbp  
 51and consist of about 80% of highly repetitive elements (repeats), is a rather complex task [1-12].  
 52The main problem of assembling such genomes is the limitations of assembler programs. As a  
 53rule, modern assemblers are designed to assemble genomes shorter or equal to the length of the

54human genome (3 Gbp). Most assemblers cannot handle the amount of input sequence data  
55required to provide the coverage needed for a high-quality assembly or take too much time and  
56computer resources. This prompts the development of new approaches in assembling large  
57genomes, including Siberian larch (*Larix sibirica* Ledeb.), which together with Siberian stone  
58pine (*Pinus sibirica* Du Tour) are the main objects of the genome project "Genomics of the key  
59boreal forest conifer species and their major phytopathogens in the Russian Federation" funded  
60by a research grant No. 14.Y26.31.0004 from the Government of the Russian Federation.

61

## 62**Methods**

### 63**A stepwise approach to assembling large genomes**

64High sequence coverage is always needed for high-quality *de novo* genome sequencing and  
65assembly. For a given average genome coverage, the coverage of individual genome regions is  
66approximately described by the Poisson distribution according to the Lander-Waterman theory  
67[13]. Insufficient coverage increases the probability of zero coverage of some genome regions.  
68Meanwhile, even a single coverage of genome regions is sufficient for their assembling using De  
69Bruijn graph based methods [14] assuming no errors and repeats.

70 To solve the problem a new stepwise approach to assembling large genomes "in parts" was  
71developed. The idea of partitioning data to perform assembly is not new. For example, in the  
72article [15] it was proposed to apply a similar two-step hierarchical approach with the aim of  
73improving the quality of assembly of bacterial genomes with very high coverage. However, the  
74approach presented in [15] does not solve the problems of assembling large and super-large  
75genomes, especially if DNA was obtained from diploid tissue.

76 In our case the assembly is also done in two steps. At the first step, the entire input pool of the  
77sequence reads is divided into several sets (parts). The size of each set is within a limit for the  
78number of reads that can be handled by the assembler program. Each set is assembled separately,

79then contigs obtained for each part are combined and used as the input data for the second step of  
80assembling.

81 With this approach in the second step of assembling, the genome coverage by the input  
82contigs no longer obeys the Poisson distribution. However, the level of coverage will not be  
83greater than the number of parts by which the original pool of reads has been partitioned, which  
84allows to bypass the limitation for the maximum amount of input data in the second step.

85 The challenge of the approach is the lower tolerance to sequencing errors and polymorphisms.  
86The ambiguity in the input sequences in the second step could lead to generating duplications in  
87the output. Therefore, the pipeline for the assembly with this approach should also include a  
88verification of the assembly for redundancy to exclude potential duplicates. We used the  
89UCLUST package [16] and self-blasting for this task.

90 It should be noted that not all assembly programs allow generating contigs with a coverage  
91below the threshold value. To overcome this obstacle in the second step of the stepwise assembly  
92either the program codes should be changed or software that does not have these limitations,  
93such as the CLC Assembly Cell (QIAGEN, Hilden, Germany) should be used. This software  
94takes into account possible sequencing errors during assembling. Thus, if there are sequencing  
95errors in the input reads, most of them will not be incorporated in the contigs generated in the  
96first step for each part of the pool. However, the problem of the stepwise assembling could be  
97insufficient coverage for each part, which can lead to shorter contigs. Since there is a restriction  
98on the minimum length of contigs in the assembling programs, such short contigs with  
99insufficient lengths will be excluded from the assembly. Therefore, to reduce the probability of  
100gaps due to excluding short contigs in the second step, one of the sets in the first step included all  
101reads from the original data pool, but to make computing possible they were used as single end  
102reads, and they were also multiplied. All steps are presented as a workflow chart in Fig. 1.

103

104**Testing of the proposed stepwise approach on the model plant species *Arabidopsis thaliana***

To test the applicability of the proposed method of stepwise assembling for *de novo* assembling of large genomes, such as in *L. sibirica* (12.03 Gbp), a genome assembly of the model plant species *Arabidopsis thaliana* obtained by the proposed method was compared with the standard *de novo* assembly of this species genome. A relatively small subset of *A. thaliana* genomic reads was selected to get a genome coverage comparable to *L. sibirica*.

As an additional argument supporting applicability of the method, the histograms of genome coverage obtained for *A. thaliana* and *L. sibirica* were compared for similarity. To construct the histograms, the genomic reads used for assembling were mapped to the assembled genomes using the bowtie software [17] for *A. thaliana* and the CLC read mapper for *L. sibirica*.

The *A. thaliana* genome contains 5 chromosomes and 135 Mbp [18]. We used the SPAdes [19], AbySS [20], CLC Assembly Cell (<https://www.qiagenbioinformatics.com/products/clc-assembly-cell>), and SOAPdenovo [21] assemblers for the traditional *de novo* assembly of the *A. thaliana* genome. The genomic paired-end reads of *A. thaliana* were downloaded from the Genbank SRA database (accession number SRR492411 [22]). The results of assembly at the level of contigs by different assemblers are presented in Fig. 2 and Additional file 1.

120

121 <Fig. 2 location>

122

The result of assembling repetitive regions of the genome depends on the number and similarity of copies of a particular type of repeat. With a small and divergent number of copies, the assembler program, as a rule, is able to separate individual copies, so that all variants of this repeat will be presented in the final contigs. With a large number of identical or nearly identical copies of the same type, it would be difficult for an assembler to separate them. The number of repeats in the genome of *A. thaliana* represents quite a significant part - according to different estimates from 23 to 32 % [23, 24]. As a result, in the final assemblies, identical repeats of the

130 same type can be represented by a single contig. This was reflected in the histogram of the contig  
 131 coverage based on the distribution of mapped reads used for assembling and presented in Fig. 3.

132

133 <Fig. 3 location>

134

135 It should also be noted that in the area of maximum coverage its distribution is more accu-

136 rately described by the corrected Poisson distribution expressed by the formula  $\frac{bL^{bx} e^{-bL}}{\Gamma(bx+1)}$ ,  
 137 where  $L$  - average coverage,  $x$  - coverage value,  $b$  - correction parameter (inversed value of ex-  
 138 tended variation) (Fig. 3, dotted line,  $b = 0.3$ ).

139 The observed coverage histogram followed the Lander-Waterman theory in general, and the  
 140 degree of coverage can be approximately described by the Poisson distribution for most of the  
 141 genome with the left side maximum peak equaled to 16 reads (Fig. 3). The exact fitting of the  
 142 coverage histogram to the Poisson distribution and the corrected (over-dispersed) Poisson  
 143 distribution was estimated using an iterative maximum likelihood-based procedure implemented  
 144 in the R statistical package. The results of these tests confirmed the fitting of the histogram to the  
 145 over-dispersed Poisson distribution around the peak value, with the reservations about semi-  
 146 qualitative description of the distribution. The left and right tails of the distribution do not obey  
 147 the provided model and should be described using another approaches. Because of this, the  
 148 goodness of the fitting depends on selection of limits around a peak value of distribution. In  
 149 reasonable limits between 0.5X and 2X of peak value, the match to over-dispersed Poisson  
 150 distribution was significant based on the Kolmogorov-Smirnov (KS) test ( $P < 0.01$ ), but the  
 151 estimated values of parameters should be anyway considered as approximate to avoid an excess  
 152 of accuracy.

153 The clearly observed “heavy tail” in the right part of the distribution for contigs with high  
 154 coverage (more than 100 reads) could be explained by the highly repetitive elements that

155 represented different parts in the original genome, but were aligned and mapped together to the  
 156 same single contigs. Therefore, the observed coverage histogram can be divided into two parts,  
 157 with a coverage less or more than 100 reads, respectively. The key observation was that the  
 158 observed coverage histogram for the *L. sibirica* genome followed the same trend that further  
 159 confirms the applicability of the proposed method (respective larch data and figures are  
 160 presented and discussed below in Results). The “heavy tails” were also observed in the coverage  
 161 histograms in metagenomics [25] and medical DNA sequencing [26].

162 The number of copies of different types of repeats in the genome is governed by different  
 163 evolutionary factors, and the simplest way to explain the heavy tail of the distribution is to use  
 164 the Zipf’s law to describe the frequencies of different types of repeats [27]. According to the  
 165 Zipf’s law, the frequencies of different types of repeats, sorted by the degree of occurrence,  
 166 should be distributed in proportion to  $1/n$ , where  $n$  is a consecutive number of the type of repeat  
 167 in the list of observed types.

168 The number of repeats with a given degree of coverage can be expressed as the derivative of  
 169 this dependence, that is, in proportion to  $1/n^2$ , where  $n$  is the degree of coverage. If the value of

170  $Z = \frac{1}{\sqrt{Y}}$  is calculated for a coverage histogram same as in Fig. 3, where  $Y$  is the percentage of  
 171 the genome with a given degree of coverage, then according to the Zipf’s law, the value of  $Z$   
 172 should directly and proportionality depend on the degree of coverage. This dependence is  
 173 demonstrated in Fig. 4 for the histogram of the observed coverage presented in Fig. 3.

174

175 <Fig. 4 location>

176

177 As it can be seen from Fig. 4, the Zipf’s law is approximately satisfied for the coverage of  
 178 more than 200 reads per site, which agrees with the abovementioned conclusion about  
 179 assembling repeats that occurred with different frequency in the genome. For a more accurate

180description of the observed dependence, it is recommended to use a distribution based on the

181Zipf-Mandelbrot law formulated as  $\frac{1}{n^k}$ , where  $k$  is generally different from unity [27].  
 182Nevertheless, the applicability of this law to genomic nucleotide sequences requires further  
 183study.

184 There are a few studies of the *A. thaliana* genome that identified different types of repeats,  
 185using, in particular, the method of clustering repeat sequences (for example, [23, 24]). According  
 186to these studies, while there was a general tendency to meet the Zipf's law for regions with a high  
 187degree coverage, individual peaks also appeared in the coverage distributions, such as in our case  
 188(Fig. 4), which can be interpreted as a manifestation of the similarity between individual types of  
 189repeats.

190 As shown in Fig. 3, the *A. thaliana* genome coverage was mostly described by a Poisson  
 191distribution with an average value of about 16 reads. To test the suggested stepwise assembling  
 192method, four sets were generated from the original pool of about 13 million reads. The first three  
 193sets included the first, second and third thirds of the original pool of reads, respectively. The  
 194fourth set also included one third of the original pool of reads, but was generated by random  
 195sampling from the mixed original pool of reads.

196 Thus, four sets of reads were generated from the original pool of reads used in the tests  
 197presented in Fig. 2 and Additional file 1. Fig. 5 and Additional file 2 presents the results of the  
 198stepwise assembly by four assemblers when each of the sets (parts) was assembled separately in  
 199the first step and then finally assembled by pooling all contigs from all four sets. It can be seen  
 200from the table that the CLC Assembly Cell demonstrated the best performance.

201

202

<Fig. 5 location>

203



Table 1 presents the results of assembly of each of the sets (parts) separately (the first step), as well as based on the pooling of contigs obtained respectively from two, three, and four sets (parts) using the CLC Assembly Cell software.

<Table 1 location>

Table 1 shows that insufficient coverage led to a significant decrease in the average contig length compared to the data in Fig. 2 and Additional file 1, but in the second step of assembling this parameter was corrected, and with the increase in the number of parts was stabilized at the level of values close to the values obtained by the different assemblers used to assemble the entire pool of reads simultaneously.

The identity of assembly obtained using parts and the stepwise method with assembly based on assembling simultaneously all reads was tested by the NUCmer software (<http://mummer.sourceforge.net>), and the highest similarity was for alignments generated by the CLC Assembly Cell (90.14%) and Abyss (95.24%) software, respectively (Fig. 5 and Additional file 2), but the former software computed the assembly with the less number of contigs and more realistic total length, and seven times faster than the latter one with the same computer hardware resources (31 vs. 217 minutes, Fig. 5 and Additional file 2).

Fig. 6 compares the genome coverage histograms for the *A. thaliana* genome assembly based on assembling the entire pool of reads simultaneously, such as in Fig. 3, and assembly based on the stepwise assembling in two steps of four parts (Table 1). It is clearly seen in Fig. 6 that the stepwise assembled genome was adequately covered by the original set of reads.

<Fig. 6 location>

The ambiguous positions in the *A. thaliana* sequencing data were estimated by aligning original *A. thaliana* reads to the assembly by Bowtie2. They represented 0.7% of genome size.

230The duplications of contigs were not detected in the final assembly, thus indicating a low level of  
231ambiguity for the assembly obtained by the suggested method.

232

### 233**The stepwise approach for the *Larix sibirica* genome assembly**

234For the assembly of the *L. sibirica* genome, four PE and three MP libraries with different insert  
235size were used (Fig. 7 and Additional file 3). At the first step, MPE libraries were decoupled and  
236used as single reads to complete a pool of reads. The pool of reads was split to four parts and  
237four sets of conigs were obtained, respectively. The CLC Assembly Cell software was selected  
238for assembling the larch genome as the best performing software.

239

240 <Fig. 7 location>

241

242 Also, a fifth set of reads was added to the analysis. This set included all reads, but the PE and  
243MPE reads were decoupled and used as single reads. This set was generated because we found  
244experimentally that the CLC Assembly Cell assembler was able to process the entire volume of  
245the *L. sibirica* sequence data, but only if the information about the length of the insertion was not  
246indicated. In this case the "Optimization of the graph using paired reads" step is skipped. In this  
247step long repeats are allowed, and scaffolding is not performed which turns out to be too much  
248computationally intense and practically prohibitive for large volume data. Therefore, this set  
249increased the representation of all reads, but they all could be used only as the single end reads at  
250this step.

251 Unlike the inbred highly homozygous plant used for the genome sequencing and assembly,  
252such as *A. thaliana*, the *L. sibirica* tree used for genome sequencing in our study represented a  
253common forest tree with a relatively high level of individual heterozygosity and, respectively,  
254high within individual biallelic variation. The number of ambiguous positions in the *L. sibirica*  
255sequencing data was estimated at level 3.0% of genome size. The presence of duplicate contigs  
256was detected in the preliminary draft assembly of *L. sibirica* obtained in the second step, thus

revealing the higher data ambiguity in the *L. sibirica* sequencing data compared to the *A. thaliana* data. To resolve the ambiguities in the second stage, the total number of all contigs resulted from the fifth set was increased by 16 folds by multiplying each contig 16 times, respectively. This trick allowed the CLC assembler to apply the majority rule when picking one of the alternative alleles, using the alleles selected in the fifth set in the first step of assembly. The same approach was used also for the *Arabidopsis thaliana* genome stepwise assembly by four different assemblers (Fig. 8 and Additional file 4). The CLC Assembly Cell again demonstrated the best performance.

265

266 &lt;Fig. 8 location&gt;

267

In addition, to verify the accuracy of the stepwise CLC Assembly Cell assembly the medium size genome (265 Mb, 2n =16) of *Prunus persica* (peach) was also assembled by both the traditional method using 24324216 sequence reads (~15X coverage) available on <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA31227> and the same stepwise approach that was used for the larch genome assembly and based on the five parts (Fig. 9 and Additional file 5). The traditional and stepwise assemblies were similar by 95.64% based on the NUCMER comparison.

275

&lt;Fig. 9 location&gt;

276

## 277 Results

### 278 Stepwise assembly of the *Larix sibirica* genome in parts

The length of the *L. sibirica* genome is about 12.03 Gb [28], and about 82 % of which consists of repeats [6-8]. The volume of the larch sequencing data obtained (11 billion paired 100 bp long reads) was hardly manageable by the available genome assemblers and more than twice the maximum amount of data that the best performing software in our test with the *Arabidopsis* data

CLC Assembly Cell can handle. Therefore, we developed a new stepwise assembly method for assembling this and other large genomes and demonstrated its consistency in computer experiments on assembling the model plant *A. thaliana* genome.

The original Siberian larch sequencing data were partitioned into five sets following mainly the procedure described for *A. thaliana* in Methods with an additional fifth set. Each set was separately assembled by the CLC Assembly Cell program. The assembly results are presented in Table 2 for each set. Only contigs with a minimum length of 200 bp were included in the final assembly.

<Table 2 location>

Total length of contigs assembled separately for each of the five sets varied from ~2.5 to ~6 Gb. The N50 parameter varied from ~300 to ~1300 bp. At the second step, individual assemblies were combined by specifying them as unpaired reads and changing the *k*-mer parameter length from 35 to 60. In addition, the mate pair (MP) reads generated from the MP libraries with 2000-10,000 bp long inserts were added to the CLC Assembly Cell input data. These reads were used at the stage of scaffolding (joining contigs into scaffolds with gaps of known expected length).

Additional scaffolding was done using BESST [29], and 228,571 additional scaffolds were generated. The scaffolding was also improved by using larch transcriptome reads and RaScaf + Bowtie2 software [30]. About 92% reads were mapped to the genome assembly and allowed us to connect 3,622 contigs into scaffolds. The assembly was finished with gap-closing using the Sealer program implemented in the last part of the Abyss pipeline [31], and 61,037 gaps were closed.

Thus, the contigs of the all five assemblies were processed and the obtained statistics presented in Table 2.

Adding the last two assemblies based on the 4 and 5 sets (Table 2) improved the final parameters (Table 3), and increased the total contig and scaffold lengths from 7.18 to 7.99 Gb

309 and from 11.04 to 12.34 Gb, respectively. The N50 parameter remained unchanged compared to  
310 the best values of partial assemblies. This is inconsistent with the results for *A. thaliana*  
311 assembly tests, but could be explained by the additional scaffolding procedure with the MP reads  
312 for the *L. sibirica* assembly.

313 <Table 3 location>

314

315 The assembly was tested for redundancy using a custom pipeline specially developed for this  
316 task, which checks for duplication taking into account possible erroneous nucleotide  
317 substitutions and indels. As a result, 74851 scaffolds were excluded. The assembly was  
318 additionally checked for vector contamination and redundancy using the UniVec database  
319 (<https://www.ncbi.nlm.nih.gov/tools/vecscreen/univec>) and the BLAST program, and as a result,  
320 10681 sequences were deleted.

321 Finally, after scaffolding, a complete Siberian larch genome of 12.34 Gb was assembled *de*  
322 *novo*. The computing time taken to assemble larch genome using 40 cores is presented in Fig. 10  
323 and Additional file 6. In total it took about 529 hours or 22 days. Therefore, the larch genome  
324 computing using the next best assembler SOAPdenovo could predictively take more than 100  
325 days.

326 <Fig. 10 location>

327

328 The histogram of the coverage for the obtained genome corresponded to the Poisson distribu-  
329 tion with extended variation in the regions with low coverage (Fig. 11A) and to the Zipf's law in  
330 the region of high coverage (Fig. 11B) and was similar to the one obtained for *A. thaliana* (Fig.  
331 13). The values for the inversed over-dispersion parameter were nearly the same for both genomes  
332 ( $0.3 \pm 0.1$ ), as it was confirmed by likelihood-based parameter estimates.

333

334 <Fig. 11 location>

335

336 The correlation presented in Fig. 11B was completely linear in the region of sufficient  
337 coverage, as expected from the Zipf's law, in contrast to the correlation for *A. thaliana*, in which  
338 many individual peaks were observed (Fig. 4). This is consistent with the results of the analysis  
339 of genomic repeats in Norway spruce [1], where it was difficult to cluster repeats and separate  
340 some types of repeats, as it can be done for many other genome sequences of eukaryotes. It also  
341 followed from our results that distribution of repeats is continuous in conifers. The presence of a  
342 large number of repeats and associated with them discontinuities in assembling can explain the  
343 smaller average contig length in comparison with the results of the *A. thaliana* genome  
344 assembling.

345 The accuracy of the stepwise CLC Assembly Cell assembly was also verified by assembling  
346 the medium size genome (265 Mb,  $2n = 16$ ) of *Prunus persica* (peach) using both methods. The  
347 assembly parameters are presented in Table 3 and the histogram of the coverage in Fig. 12. Both  
348 observed and expected distributions of the peach genome coverage were similar to those for  
349 *Arabidopsis* (Figs 2, 3, 5) and Siberian larch (Fig. 11) genomes.

350

351

&lt;Fig. 12 location&gt;

352

353 The negative binomial distribution or the over-dispersed Poisson distribution is often used to  
354 describe genome coverage histograms, but, to our best knowledge, the effect of overdispersion  
355 was not systematically studied in the context of genome assemblies (but see [25, 32]). However,  
356 the similar values of the over-dispersion parameter for three assembled genomes confirmed by  
357 the KS tests could serve as an additional argument that the proposed method could be adequately  
358 scaled to the assembly of large genomes.

359

## 360 Discussion

361 Testing of the proposed stepwise approach for assembling genomes in parts on the model plant  
362 species *A. thaliana* showed that despite some deterioration of the distribution parameters of the  
363 contig lengths in the final assembly compared to normal assembling using the CLC Assembly  
364 Cell, the result of the stepwise assembling was comparable with the results of assembling all data  
365 simultaneously using different assemblers. Comparison of the lengths of the obtained genomes  
366 and histograms of the coverage obtained by different methods also allows us to state that the  
367 stepwise assembling by parts allows obtaining a consistent and reliable genome assembly  
368 corresponding to the original biological material.

369 The analysis of the coverage histograms carried out for *A. thaliana*, *Prunus persica* (peach)  
370 and larch showed a tendency to satisfy the Zipf's law for the frequency of repeats and provided  
371 additional grounds for concluding that the stepwise assembly approach by parts is applicable for  
372 assembling large genomes, such as the Siberian larch genome. The interpretation of the coverage  
373 histograms using the Zipf's law made it possible also to clarify the idea of the statistical  
374 regularities characterizing the evolutionary mechanisms of multiplication of repeats in different  
375 plant species.

376

## 377 Conclusion

378 Using the new stepwise *de novo* assembling method presented in the paper the genome of  
379 Siberian larch, *Larix sibirica* Ledeb. (12.34 Gbp) was for the first time completely assembled *de*  
380 *novo* by the CLC Assembly Cell assembler. It is the first genome assembly for any larch species  
381 in addition to only five other conifer genomes sequenced and assembled for *Picea abies* [1],  
382 *Picea glauca* [2], *Pinus taeda* [3-5, 9, 11], *Pinus lambertiana* [10] and *Pseudotsuga menziesii*  
383 var. *menziesii* [12]. Presented approach makes assembling feasible for very large genomes with a  
384 reasonable computing time and without engaging huge computing resources. The assemblies  
385 produced by this approach are still of reasonable quality allowing their annotation and further  
386 use.

387

## 388 **Declarations**

389

## 390 **Abbreviations**

391 Gb: Giga Base; bp: base pair; HPC: High Performance Computing

392

## 393 **Acknowledgements**

394 We thank the Department of High Performance Computing for their help with computing using  
395 their HPC cluster at the Siberian Federal University.

396

## 397 **Funding**

398 This study was funded by a research grant No. 14.Y26.31.0004 from the Government of the  
399 Russian Federation. No funding agency played any role in the design or conclusion of this study.  
400 Publication costs are funded by the BioMed Central Membership of the University of Göttingen.

401

## 402 **Availability of data and materials**

403 This manuscript describes published software and a new developed pipeline (source code is  
404 available from the authors on request). The sequence reads and obtained scaffolds are publicly  
405 available under the NCBI Genbank BioProject accession number PRJNA393226.

406

## 407 **Authors' contributions**

408 KVK & DAK conceived and developed original idea. SIF, VVS, ANC, SVM & YAP wrote the  
409 codes, developed the pipeline and implemented parallelization. YAP and NVO processed original  
410 sequencing reads and generated the original data. DAK, VVS, SIF & KVK together wrote the  
411 first draft manuscript. All authors read, revised and approved the final manuscript.

412



### 413Competing interests

414The authors declare that they have no competing interests.

415

### 416Consent for publication

417Not applicable.

418

### 419Ethics approval and consent to participate

420Not applicable.

421

### 422References

4231. Nystedt B, Street NR, Wetterbom A, Zuccolo A, Lin YC, Scofield DG, Vezzi F, Delhomme  
424 N, Giacomello S, Alexeyenko A, Vicedomini R, Sahlin K, Sherwood E, Elfstrand M, Gram-  
425 zow L, Holmberg K, Hällman J, Keech O, Klasson L, Koriabine M, Kucukoglu M, Käller  
426 M, Luthman J, Lysholm F, Niittylä T, Olson A, Rilakovic N, Ritland C, Rosselló JA, Sena J,  
427 Svensson T, Talavera-López C, Theißen G, Tuominen H, Vanneste K, Wu ZQ, Zhang B,  
428 Zerbe P, Arvestad L, Bhalerao R, Bohlmann J, Bousquet J, Garcia Gil R, Hvidsten TR, de  
429 Jong P, MacKay J, Morgante M, Ritland K, Sundberg B, Thompson SL, Van de Peer Y, An-  
430 dersson B, Nilsson O, Ingvarsson PK, Lundeberg J, Jansson S. Norway spruce genome se-  
431 quence and conifer genome evolution. *Nature*. 2013;497:579-584. doi:10.1038/nature12211
4322. Birol I, Raymond A, Jackman SD, Pleasance S, Coope R, Taylor GA, Saint Yuen MM, Keel-  
433 ing CI, Brand D, Vandervalk BP, Kirk H, Pandoh P, Moore RA, Zhao YJ, Mungall AJ,  
434 Jaquish B, Yanchuk A, Ritland C, Boyle B, Bousquet J, Ritland K, MacKay J, Bohlmann J,  
435 Jones SJM. Assembling the 20 Gb white spruce (*Picea glauca*) genome from whole-genome  
436 shotgun sequencing data. *Bioinformatics*. 2013;29(12):1492-1497. doi:10.1093/bioinformat-  
437 ics/btt178

4383. Neale DB, Wegrzyn JL, Stevens KA, Zimin AV, Puiu D, Crepeau MW, Cardeno C,  
 439 Koriabine M, Holtz-Morris AE, Liechty JD, Martinez-Garcia PJ, Vasquez-Gross HA, Lin  
 440 BY, Zieve JJ, Dougherty WM, Fuentes-Soriano S, Wu LS, Gilbert D, Marcais G, Roberts M,  
 441 Holt C, Yandell M, Davis JM, Smith KE, Dean JF, LorenzWW, Whetten RW, Sederoff R,  
 442 Wheeler N, McGuire PE, Main D, Loopstra CA, Mockaitis K, deJong PJ, Yorke JA, Salzberg  
 443 SL, Langley CH. Decoding the massive genome of loblolly pine using haploid DNA and  
 444 novel assembly strategies. *Genome Biol.* 2014;15(3):R59. doi:10.1186/gb-2014-15-3-r59
4454. Wegrzyn JL, Liechty JD, Stevens KA, Wu LS, Loopstra CA, Vasquez-Gross HA, Dougherty  
 446 WM, Lin BY, Zieve JJ, Martinez-Garcia PJ, Holt C, Yandell M, Zimin AV, Yorke JA,  
 447 Crepeau MW, Puiu D, Salzberg SL, Dejong PJ, Mockaitis K, Main D, Langley CH, Neale  
 448 DB. Unique features of the loblolly pine (*Pinus taeda* L.) megagenome revealed through  
 449 sequence annotation. *Genetics.* 2014;196:891–909. doi:10.1534/genetics.113.159996
4505. Zimin A, Stevens KA, Crepeau MW, Holtz-Morris A, Koriabine M, Marcais G. Puiu D,  
 451 Roberts M, Wegrzyn J, de Jong P, Neale D, Salzberg S, Yorke J, Langley C. Sequencing and  
 452 Assembly of the 22-Gb Loblolly Pine Genome. *Genetics.* 2014;196(3):875-890.  
 453 doi:10.1534/genetics.113.159715
4546. Krutovsky KV, Oreshkova NV, Putintseva YA, Ibe AA, Deich KO, Shilkina EA. Preliminary  
 455 results of de novo whole genome sequencing of Siberian larch (*Larix sibirica* Ledeb.) and  
 456 Siberian stone pine (*Pinus sibirica* Du Tour.). *Siberian Journal of Forest Science.*  
 457 2014;1(4):79-83.
4587. Oreshkova NV, Putintseva YuA, Kuzmin DA, Sharov VV, Biryukov VV, Makolov SV, De-  
 459 ich KO, Ibe AA, Shilkina EA, Krutovsky KV. Genome sequencing and assembly of Siberian  
 460 larch (*Larix sibirica* Ledeb.) and Siberian pine (*Pinus sibirica* Du Tour) and preliminary  
 461 transcriptome data // *Proceedings of the 4th International Conference on Conservation of*  
 462 *Forest Genetic Resources in Siberia.* Barnaul, Russia, 24–29 August, 2015; p. 127–128.
4638. Krutovsky KV, Putintseva YuA, Oreshkova NV, Kuzmin DA, Pavlov IN, Sharov VV,  
 464 Biryukov VV, Makolov SV, Deych KO, Bondar EI, Ushakova OA, Ibe AA, Shilkina EA,

- 465 Sadovsky MG, Vaganov EA. *Pinus sibirica* and *Larix sibirica* whole genome *de novo* se-  
466 quencing. IUFRO Genomics and Forest Tree Genetics Conference, May 30 - June 3, 2016,  
467 Arcachon, France. Oral presentation. Book of Abstracts. 2016; p. 39  
468 (<https://colloque.inra.fr/iufro2016/Programme>).
4699. Sadovsky MG, Putintseva YA, Birukov VV, S. Novikova, Krutovsky KV. *De novo* assembly  
470 and cluster analysis of Siberian larch transcriptome and genome. Lecture Notes in  
471 Bioinformatics. 2016;9656:455-464. doi:10.1007/978-3-319-31744-1 41.
47210. Stevens KA, Wegrzyn JL, Zimin A, Puiu D, Crepeau M, C Cardeno, R Paul, D Gonzalez-  
473 Ibeas, M Koriabine, AE Holtz-Morris, PJ Martínez-García, UU Sezen, G Marçais, K  
474 Jermstad, PE McGuire, CA Loopstra, JM Davis, A Eckert, P de Jong, JA Yorke, SL Salzberg,  
475 DB Neale, Langley CH. Sequence of the sugar pine megagenome. Genetics.  
476 2016;204(4):1613–1626. doi:10.1534/genetics.116.193227.
47711. Zimin A, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, Langley CH, Neale  
478 DB, Salzberg SL. An improved assembly of the loblolly pine mega-genome using long-read  
479 single-molecule sequencing. Gigascience. 2017;6:1–4. doi:10.1093/gigascience/giw016.
48012. Neale DB, McGuire PE, Wheeler NC, Stevens KA, Crepeau MW, Cardeno C, Zimin AV,  
481 Puiu D, Pertea GM, Sezen UU, Casola C, Koralewski TE, Paul R, Gonzalez-Ibeas D, Zaman  
482 S, Cronn R, Yandell M, Holt C, Langley CH, Yorke JA, Salzberg SL, Wegrzyn JL. The Dou-  
483 glas-Fir Genome Sequence Reveals Specialization of the Photosynthetic Apparatus in  
484 Pinaceae. 2017. G3: Genes, Genomes, Genetics. 2017;7(9):3157-3167.  
485 doi:10.1534/g3.117.300078.
48613. Lander ES, Waterman MS. Genomic mapping by fingerprinting random clones: a mathemat-  
487 ical analysis. Genomics.1988 2(3):231-239.
48814. Compeau PEC, Pevzner PA, Tesler G. How to apply de Bruijn graphs to genome assembly.  
489 Nat Biotechnology. 2011;29(11):987–991.

49015. Al-Okaily AA. HGA: *de novo* genome assembly method for bacterial genomes using high  
491 coverage short sequencing reads. BMC Genomics. 2016; 17:193.  
492 <https://doi.org/10.1186/s12864-016-2515-7>.
49316. Edgar RC. Search and clustering orders of magnitude faster than BLAST. Bioinformatics.  
494 2010;26(19):2460–2461.
49517. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of  
496 short DNA sequences to the human genome. Genome Biol. 2009;10:R25.
49718. Bennett MD, Leitch IJ, Price HJ, Johnston JS. Comparisons with *Caenorhabditis* (~100 Mb)  
498 and *Drosophila* (~175 Mb) using flow cytometry show genome size in *Arabidopsis* to be  
499 ~157 Mb and thus ~25 % larger than the Arabidopsis Genome Initiative estimate of ~125  
500 Mb. Annals of Botany. 2003;91:547-557.
50119. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,  
502 Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,  
503 Alekseyev MA, Pevzner PA. SPAdes: a new genome assembly algorithm and its applications  
504 to single-cell sequencing. J Comput Biol. 2012;19(5):455-477.
50520. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I. ABySS: a parallel  
506 assembler for short read sequence data. Genome Res. 2009; 19:1117-1123.
50721. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, Tang J, Wu G,  
508 Zhang H, Shi Y, Liu Y, Yu C, Wang B, Lu Y, Han C, Cheung DW, Yiu SM, Peng S, Xiaoqian  
509 Z, Liu G, Liao X, Li Y, Yang H, Wang J, Lam TW, Wang J. SOAPdenovo2: an empirically  
510 improved memory-efficient short-read *de novo* assembler. Gigascience. 2012;1:18.
51122. Schneeberger K, Ossowski S, Ott F, Klein JD, Wang X, Lanz C, Smith LM, Cao J, Fitz J,  
512 Warthmann N, Henz SR, Huson DH, Weigel D. Reference-guided assembly of four diverse  
513 *Arabidopsis thaliana* genomes. Proc Natl Acad Sci USA. 2011;108(25):10249-10254.

51423. Maumus F, Quesneville H. Ancestral repeats have shaped epigenome and genome  
515 composition for millions of years in *Arabidopsis thaliana*. Nat Communications.  
516 2014;5:4104.
51724. Maumus F, Quesneville H. Deep Investigation of *Arabidopsis thaliana* Junk DNA Reveals a  
518 Continuum between Repetitive Elements and Genomic Dark Matter. PLOS ONE.  
519 2014;9(4):e94101.
52025. Lindner MS, Kollock M, Zickmann F, Renard BY. Analyzing genome coverage profiles with  
521 applications to quality control in metagenomics, Bioinformatics. 2013;29(10):1260–1267.
52226. Wendl MC, Wilson RK. Aspects of coverage in medical DNA sequencing. BMC  
523 Bioinformatics. 2008;9:239. doi:10.1186/1471-2105-9-239.
52427. Mantegna RN, Buldyrev SV, Goldberger AL, Havlin S, Peng CK, Simons M, Stanley HE.  
525 Linguistic Features of Noncoding DNA Sequences. Phys Rev Lett. 1994;73(23):3169–3172.
52628. Ohri D, Khoshoo TN. Genome size in gymnosperms. Plant Systematics and Evolution.  
527 1986;153:119-132.
52829. Sahlin K, Vezzi F, Nystedt B, Lundeberg J, Arvestad L. BESST - efficient scaffolding of  
529 large fragmented assemblies. BMC Bioinformatics. 2014;15(1):281. doi:10.1186/1471-  
530 2105-15-281.
53130. Song L, Shankar DS, Florea L. Rascaf: Improving Genome Assembly with RNA Sequencing  
532 Data. Plant Genome. 2016;9(3). doi: 10.3835/plantgenome2016.03.0027.
53331. Paulino D, Warren RL, Vandervalk BP, Raymond A, Jackman SD, Birol I. Sealer: a scalable  
534 gap-closing application for finishing draft genomes. BMC Bioinformatics. 2015;16(1):230.  
535 doi:10.1186/s12859-015-0663-4.
53632. Oberg AL, Bot BM, Grill DE, Poland GA, Therneau TM. Technical and biological variance  
537 structure in mRNA-Seq data: life in the real world. BMC Genomics. 2012; 13:304.  
538 doi:10.1186/1471-2164-13-304.
- 539

**Table 1** Results of the *Arabidopsis thaliana* genome stepwise assembling in four sets (parts) using the CLC Assembly Cell software

Assembly part	Total length, bp	Contigs		
		N50, bp	number	mean length, bp
1 <sup>a</sup>	101200000	1586	110067	919
2 <sup>a</sup>	101200000	1601	109903	920
3 <sup>a</sup>	101.2	1595	110119	919
4 <sup>b</sup>	101.2	1586	110384	916
1+2	113.2	3225	64543	1753
1+2+3	116.6	3861	60606	1923
1+2+3+4	113.7	4325	52576	2161

<sup>a</sup>Represents approximately 1/3 of all original reads; <sup>b</sup>Represents also approximately 1/3 of all original reads, but randomly selected

**Table 2** The assembly results of the five sets generated from the original *Larix sibirica* genome sequencing data

Set	Number of contigs	N50, bp	Maximum length, bp	Total length, Gbp
1	7870837	310	56157	2.566
2	5469129	535	65362	2.549
3	5449065	1383	157662	4.319
4	4677717	1092	91349	3.117
5	13244672	475	46203	5.937

**Table 3** The final stepwise *Larix sibirica* genome assembly based on five sets and the MP reads

Assembly*	Number, mln	N50, bp	Maximum length, bp	Total length, Gbp
Contigs	12.40	1074	128642	7.99
Scaffolds	11.33	6443	354326	12.34

\*Minimum contig length used for assembling was 200 bp.

## 547 **Additional files**

548 **Additional file 1:** Table S1: The results of the traditional *de novo Arabidopsis thaliana* genome  
549 assembly generated by four different assemblers

550 **Additional file 2:** Table S2: Results of the *Arabidopsis thaliana* genome stepwise assembly by  
551 different assemblers using raw reads partitioned into four sets

552 **Additional file 3:** Table: S3 Sequencing libraries and generated sequence data used for the *Larix*  
553 *sibirica* genome assembly

554 **Additional file 4:** Table S4 Results of the *Arabidopsis thaliana* genome stepwise assembly by  
555 four different assemblers using raw reads partitioned into five sets following approach used for  
556 assembling of the *Larix sibirica* genome

557 **Additional file 5:** Table S5 The traditional and stepwise CLC Assembly Cell genome assembly  
558 parameters for peach (*Prunus persica*)

559 **Additional file 6:** Table S6 The computing time taken to assemble each set and the complete  
560 *Larix sibirica* genome using 40 cores

561

562

## 563 **Figure titles (max 15 words) and legends (max 300 words)**

564

565 **Fig. 1** Stepwise assembly workflow chart

566

567 **Fig. 2** The results of the traditional *de novo Arabidopsis thaliana* genome assembly generated by  
568 four different assemblers. Minimum contig length used for assembling was 200 bp

569

570 **Fig. 3** Histogram of the *Arabidopsis thaliana* genome coverage by the mapped reads used for the  
571 genome assembly generated by the CLC Assembly Cell software (solid line). Expected and  
572 corrected Poisson distributions are represented by dashed and dotted lines, respectively. The  
573 number of reads (degree of the genome coverage) is on the horizontal axis; the logarithmic  
574 proportion of the genome with such degree of coverage is on the vertical axis

575

576 **Fig. 4** Dependence of the transformed value of the fraction of the genome coverage  $Z$  on the  
577 level of coverage. Solid line represents linear dependency calculated by the least square fit

578

**Fig. 5** Results of the *Arabidopsis thaliana* genome stepwise assembly by different assemblers using raw reads partitioned into four sets. Minimum contig length used for assembling was 200 bp

**Fig. 6** Comparison of the *Arabidopsis thaliana* genome coverage histograms obtained for the genome assembly assembled by the CLC Assembly Cell using all reads simultaneously (solid line) and the stepwise method with two steps and four parts (dotted line)

**Fig. 7** Sequence coverage for seven sequencing libraries used for the *Larix sibirica* genome assembly

**Fig. 8** Results of the *Arabidopsis thaliana* genome stepwise assembly by four different assemblers using raw reads partitioned into five sets following approach used for assembling of the *Larix sibirica* genome. Minimum contig length used for assembling was 200 bp

**Fig. 9** The traditional and stepwise CLC Assembly Cell genome assembly parameters for peach (*Prunus persica*). Minimum contig length used for assembling was 200 bp.

**Fig. 10** The computing time (number of hours) taken to assemble each set and the complete *Larix sibirica* genome using 40 cores

**Fig. 11 A:** Observed distribution of Siberian larch genome coverage (solid line) and expected from the corrected Poisson distribution (dotted line) with average coverage value equalled 7 and correction parameter  $b=0.3$ . **B:** Dependence of the transformed degree of genome coverage  $Z$  on the Siberian larch genome coverage (solid line). The dashed line represents linear dependency calculated by the least square fit and fully coincides with the solid line.

**Fig. 12 A:** Observed distribution of *Prunus persica* (peach) genome coverage (solid line) and expected from the corrected Poisson distribution (dotted line) with average coverage value equalled 15 and correction parameter  $b=0.3$ . **B:** Dependence of the transformed degree of genome coverage  $Z$  on the peach genome coverage (solid line). The dashed line represents linear dependency calculated by the least square fit and fully coincides with the solid line.